

# CS 6784 Paper Presentation

## Conditional Random Fields: Probabilistic Models for Segmenting and Labeling Sequence Data

*John Lafferty, Andrew McCallum, Fernando C. N. Pereira*

Presenters: Brad Gulko and Stephanie Hyland

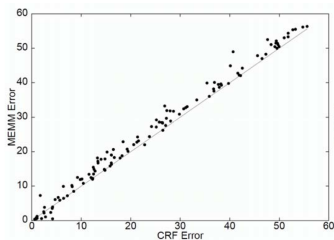
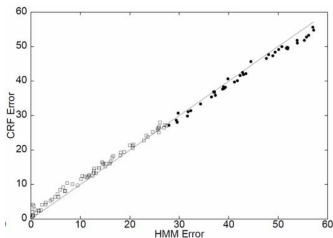
February 20, 2014

## Main Contribution Summary

- This 2001 paper introduced the **Conditional Random Field** (CRF).
- Describes efficient representation of field potentials in terms of features.
- Provides two algorithms for finding Maximum Likelihood parameter values.
- Provides some really unconvincing examples...

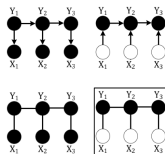
# Main Contribution Summary

... examples are NOT the strongest point of this paper



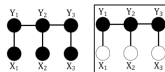
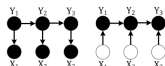
# Talk Structure

- Brad
  - CRF in context



# Talk Structure

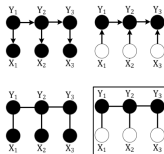
- Brad
  - CRF in context
  - The Label Bias Problem



??  
RIB  $\equiv$  ROB

# Talk Structure

- Brad
  - CRF in context
  - The Label Bias Problem
- Stephanie
  - Parameter Estimation



??  
RIB  $\stackrel{??}{=} \text{ROB}$

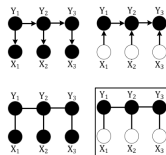
$$p(y) = \frac{1}{Z(\alpha)} \exp \left( \sum_{\text{edge}, s} \lambda_s \psi_s(y, y_{i-1}, x) + \sum_{\text{node}, s} \mu_s \phi_s(y, y_{i-1}, x) \right)$$

$$\theta = (\lambda_1, \lambda_2, \dots; \mu_1, \mu_2, \dots)$$

$$\theta_{t+1} = \theta_t - [Hf(\theta_t)]^{-1} \nabla f(\theta_t)$$

# Talk Structure

- Brad
  - CRF in context
  - The Label Bias Problem
- Stephanie
  - Parameter Estimation
  - Experiments
  - Conclusion

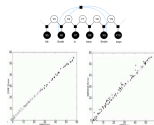


??  
RIB  $\neq$  ROB

$$p(y|x) = \frac{1}{Z(x)} \exp \left( \sum_{i \in \mathcal{V}, k} \lambda_i \phi_i(v, y_i, x) + \sum_{i \in \mathcal{V}, k} \mu_{i,k} \psi_{i,k}(v, y_i, x) \right)$$

$$\theta = (\lambda_1, \lambda_2, \dots; \mu_1, \mu_2, \dots)$$

$$\theta_{t+1} = \theta_t - [Hf(\theta_t)]^{-1} \nabla f(\theta_t)$$



## CRF in Context

- $\mathbf{X} = \{X_1, X_2, \dots\}$  be a set of observed RV
- $\mathbf{Y} = \{Y_1, Y_2, \dots\}$  be a set of label RV
- $X, Y$  be a set of joint observations of  $\mathbf{X}, \mathbf{Y}$

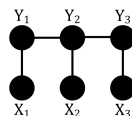
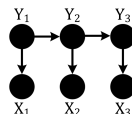
	Generative	Discriminative
Directed	???	???
Undirected	???	???



## CRF in Context

- $\mathbf{X} = \{X_1, X_2, \dots\}$  be a set of observed RV
- $\mathbf{Y} = \{Y_1, Y_2, \dots\}$  be a set of label RV
- $X, Y$  be a set of joint observations of  $\mathbf{X}, \mathbf{Y}$

	Generative	Discriminative
Directed	HMM	???
Undirected	MRF	???

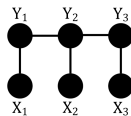
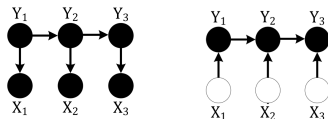


## CRF in Context

- $\mathbf{X} = \{X_1, X_2, \dots\}$  be a set of observed RV
- $\mathbf{Y} = \{Y_1, Y_2, \dots\}$  be a set of label RV
- $X, Y$  be a set of joint observations of  $\mathbf{X}, \mathbf{Y}$

	Generative	Discriminative
Directed	HMM	ME-MM
Undirected	MRF	???

In 2001, HMM, ME-MM and MRF were well known,

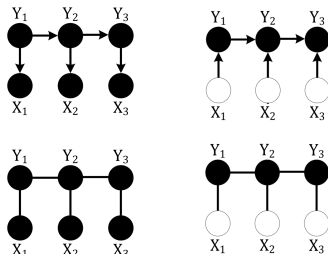


## CRF in Context

- $\mathbf{X} = \{X_1, X_2, \dots\}$  be a set of observed RV
- $\mathbf{Y} = \{Y_1, Y_2, \dots\}$  be a set of label RV
- $X, Y$  be a set of joint observations of  $\mathbf{X}, \mathbf{Y}$

	Generative	Discriminative
Directed	HMM	ME-MM
Undirected	MRF	<b>CRF</b>

In 2001, HMM, ME-MM and MRF were well known, the paper presents the CRF.



# Generative vs. Discriminative

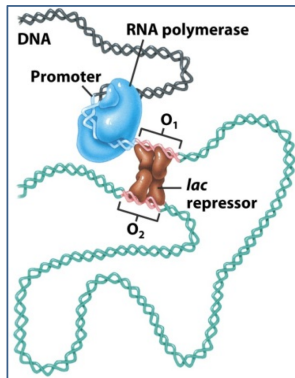
- Generative: maximise joint  $P(Y, X) = P(Y|X)P(X)$
- Discriminative: maximise conditional  $P(Y|X)$
- When is Discriminative helpful?
  - Tractability requires independence

# Generative vs. Discriminative

- Generative: maximise joint  $P(Y, X) = P(Y|X)P(X)$
- Discriminative: maximise conditional  $P(Y|X)$
- When is Discriminative helpful?
  - Tractability requires independence
  - ...but sometimes there are important correlations in  $X$ .

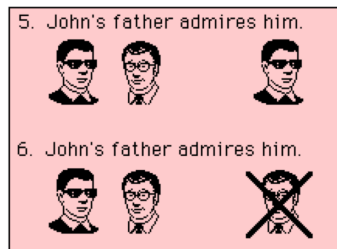
## Examples: important correlations

- Long range interactions in human genomics



## Examples: important correlations

- Long range interactions in human genomics
- Pronoun definition and binding



## Examples: important correlations

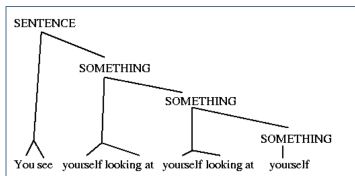
- Long range interactions in human genomics
- Pronoun definition and binding
- Context in whole scene image recognition





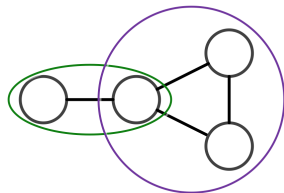
## Examples: important correlations

- Long range interactions in human genomics
- Pronoun definition and binding
- Context in whole scene image recognition
- Recursive structure in language



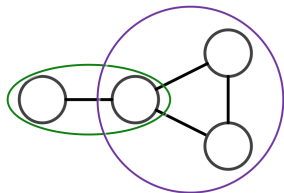
## Directed vs. Undirected

- For a graphical model  $\mathbf{G}(\mathbf{E}, \mathbf{V})$  with joint potential  $\Psi(\mathbf{V})$ .
- Let  $\mathbf{C}$  be the set of **cliques** (fully connected subgroups) in  $\mathbf{G}$ , with  $c \in \mathbf{C}$  having edges  $\mathbf{E}_c$  and vertices  $\mathbf{V}_c$ .



## Directed vs. Undirected

- For a graphical model  $\mathbf{G}(\mathbf{E}, \mathbf{V})$  with joint potential  $\Psi(\mathbf{V})$ .
- Let  $\mathbf{C}$  be the set of **cliques** (fully connected subgroups) in  $\mathbf{G}$ , with  $c \in \mathbf{C}$  having edges  $\mathbf{E}_c$  and vertices  $\mathbf{V}_c$ .
- Finally,  $Dom(\mathbf{V})$  is the set of all values assumable by the random variables,  $\mathbf{V}(= \mathbf{X} \cup \mathbf{Y})$ .



$$P(\mathbf{V}) = \frac{1}{Z} \Psi(\mathbf{V}), \quad Z = \sum_{v \in Dom(\mathbf{V})} \Psi(v)$$

## Directed vs. Undirected, continued

- Compactness requires factorization (Hammersley-Clifford, 1971):

$$\Psi(\mathbf{v}) = \prod_{c \in \mathcal{C}} \psi_c(\mathbf{v}_c)$$

## Directed vs. Undirected, continued

- Compactness requires factorization (Hammersley-Clifford, 1971):

$$\Psi(\mathbf{v}) = \prod_{c \in \mathcal{C}} \psi_c(\mathbf{v}_c)$$

- Directed: local Normalization -

$$\forall c \in \mathcal{C}, \quad \sum_{v \in \text{Dom}(\mathbf{v}_c)} \psi_c(v) = 1$$

## Directed vs. Undirected, continued

- Compactness requires factorization (Hammersley-Clifford, 1971):

$$\Psi(\mathbf{v}) = \prod_{c \in \mathcal{C}} \Psi_c(\mathbf{v}_c)$$

- Directed: local Normalization - each  $\Psi_c$  is a *probability*.

$$\forall c \in \mathcal{C}, \quad \sum_{v \in \text{Dom}(\mathbf{v}_c)} \Psi_c(v) = 1$$

## Directed vs. Undirected, continued

- Compactness requires factorization (Hammersley-Clifford, 1971):

$$\Psi(\mathbf{v}) = \prod_{c \in \mathcal{C}} \Psi_c(\mathbf{v}_c)$$

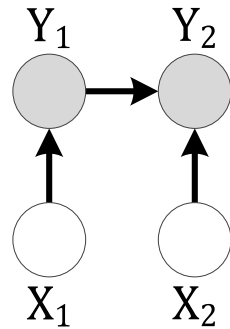
- Directed: local Normalization - each  $\Psi_c$  is a *probability*.

$$\forall c \in \mathcal{C}, \quad \sum_{v \in \text{Dom}(\mathbf{v}_c)} \Psi_c(v) = 1$$

- Undirected: Global Normalization - relaxes this constraint...  
but what does it buy us?

## The Label Bias Problem: Conditional Markov Model (EM-MM)

Toy Problem – fragment of a ME-MM



$$Y_1 \in \{1,2\} \quad Y_2 \in \{3,4\}$$

Training Data:

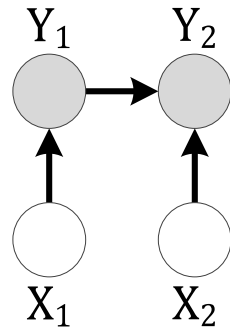
$$8: X = \{RI\} \quad Y = \{13\}$$

$$2: X = \{RO\} \quad Y = \{24\}$$



# The Label Bias Problem: Conditional Markov Model (EM-MM)

Toy Problem – fragment of a ME-MM



$$Y_1 \in \{1,2\} \quad Y_2 \in \{3,4\}$$

Training Data:

$$8: X = \{RI\} \quad Y = \{13\}$$

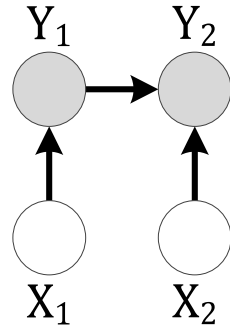
$$2: X = \{RO\} \quad Y = \{24\}$$

$P(Y_1 X_1)$	$Y_1$	
	$I$	$2$
$X_1=R$	0.8	0.2

Rel. Joint $\Psi(Y_2, X_2, Y_1)$		$Y_2$	
		$3$	$4$
$Y_1=1$	$X_2=I$	8	$\epsilon$
	$X_2=O$	$\epsilon$	$\epsilon$
$Y_1=2$	$X_2=I$	$\epsilon$	$\epsilon$
	$X_2=O$	$\epsilon$	2

# The Label Bias Problem: Conditional Markov Model (EM-MM)

Toy Problem – fragment of a ME-MM



$$Y_1 \in \{1,2\} \quad Y_2 \in \{3,4\}$$

Training Data:

$$8: X = \{RI\} \quad Y = \{13\}$$

$$2: X = \{RO\} \quad Y = \{24\}$$

$P(Y_1 X_1)$	$Y_1$	
	$I$	$2$
$X_1=R$	0.8	0.2

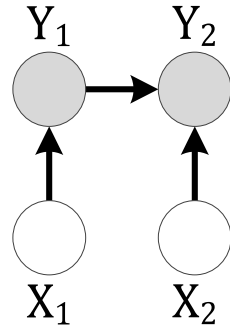
Rel. Joint $\Psi(Y_2, X_2, Y_1)$		$Y_2$	
		$3$	$4$
$Y_1=1$	$X_2=I$	8	$\epsilon$
	$X_2=O$	$\epsilon$	$\epsilon$
$Y_1=2$	$X_2=I$	$\epsilon$	$\epsilon$
	$X_2=O$	$\epsilon$	2



Conditional $P(Y_2 X_2, Y_1)$		$Y_2$	
		$3$	$4$
$Y_1=1$	$X_2=I$	$1-\epsilon$	$\epsilon$
	$X_2=O$	0.5	0.5
$Y_1=2$	$X_2=I$	0.5	0.5
	$X_2=O$	$\epsilon$	$1-\epsilon$

# The Label Bias Problem: Conditional Markov Model (EM-MM)

Toy Problem – fragment of a ME-MM



$$Y_1 \in \{1,2\} \quad Y_2 \in \{3,4\}$$

Training Data:

$$8: X = \{RI\} \quad Y = \{13\}$$

$$2: X = \{RO\} \quad Y = \{24\}$$

$P(Y_1 X_1)$	$Y_1$	
	I	2
$X_1=R$	0.8	0.2

Rel. Joint $\Psi(Y_2, X_2, Y_1)$		$Y_2$	
		3	4
$Y_1=1$	$X_2=I$	8	$\epsilon$
	$X_2=O$	$\epsilon$	$\epsilon$
$Y_1=2$	$X_2=I$	$\epsilon$	$\epsilon$
	$X_2=O$	$\epsilon$	2

Conditional $P(Y_2 X_2, Y_1)$		$Y_2$	
		3	4
$Y_1=1$	$X_2=I$	$1-\epsilon$	$\epsilon$
	$X_2=O$	0.5	0.5
$Y_1=2$	$X_2=I$	0.5	0.5
	$X_2=O$	$\epsilon$	$1-\epsilon$

Viterbi is  $P(Y_1, Y_2|X)$   
 $= P(Y_2|Y_1, X)P(Y_1|X)$   
 $= P(Y_1|X_1)P(Y_2|Y_1, X_2)$

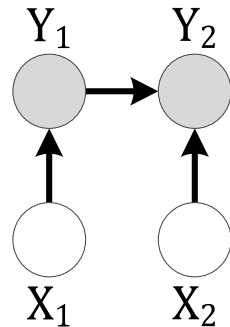
Lets try it for  $X = \{RO\}$

$Y_1, Y_2$	$P(Y_1 R)$	$P(Y_2 Y_1, O)$	$P(Y_1, Y_2 RO)$
1,3	0.8	0.5	
1,4	0.8	0.5	
2,3	0.2	$\epsilon$	
2,4	0.2	$1-\epsilon$	

Which labeling wins?

# The Label Bias Problem: Conditional Markov Model (EM-MM)

Toy Problem – fragment of a ME-MM



$$Y_1 \in \{1,2\} \quad Y_2 \in \{3,4\}$$

Training Data:

$$8: X = \{RI\} \quad Y = \{13\}$$

$$2: X = \{RO\} \quad Y = \{24\}$$

$P(Y_1 X_1)$	$Y_1$	
	1	2
$X_1=R$	0.8	0.2

Rel. Joint $\Psi(Y_2, X_2, Y_1)$		$Y_2$	
		3	4
$Y_1=1$	$X_2=I$	8	$\epsilon$
	$X_2=O$	$\epsilon$	$\epsilon$
$Y_1=2$	$X_2=I$	$\epsilon$	$\epsilon$
	$X_2=O$	$\epsilon$	2

Conditional $P(Y_2 X_2, Y_1)$		$Y_2$	
		3	4
$Y_1=1$	$X_2=I$	$1-\epsilon$	$\epsilon$
	$X_2=O$	0.5	0.5
$Y_1=2$	$X_2=I$	0.5	0.5
	$X_2=O$	$\epsilon$	$1-\epsilon$

$$\begin{aligned} \text{Viterbi is } & P(Y_1, Y_2|X) \\ &= P(Y_2|Y_1, X)P(Y_1|X) \\ &= P(Y_1|X_1)P(Y_2|Y_1, X_2) \end{aligned}$$

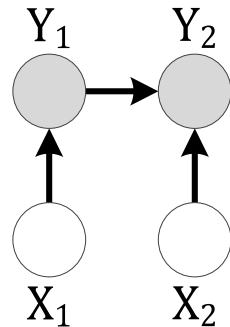
Lets try it for  $X = \{RO\}$

$Y_1, Y_2$	$P(Y_1 R)$	$P(Y_2 Y_1, O)$	$P(Y_1, Y_2 RO)$
1,3	0.8	0.5	0.4
1,4	0.8	0.5	0.4
2,3	0.2	$\epsilon$	$\epsilon$
2,4	0.2	$1-\epsilon$	0.2

But we want  
 $Y = \{2,4\}$   
What happened?

# The Label Bias Problem: Conditional Markov Model (EM-MM)

Toy Problem – fragment of a ME-MM



$$Y_1 \in \{1,2\} \quad Y_2 \in \{3,4\}$$

Training Data:

$$8: X = \{RI\} \quad Y = \{13\}$$

$$2: X = \{RO\} \quad Y = \{24\}$$

Local Normalization  
requires a  
probability.... So..

$$\frac{\epsilon}{2\epsilon} \Rightarrow \frac{1}{2}$$

Rel. Joint $\Psi(Y_2, X_2, Y_1)$		$Y_2$	
		3	4
$Y_1=1$	$X_2=I$	8	$\epsilon$
	$X_2=O$	$\epsilon$	$\epsilon$
$Y_1=2$	$X_2=I$	$\epsilon$	$\epsilon$
	$X_2=O$	$\epsilon$	2

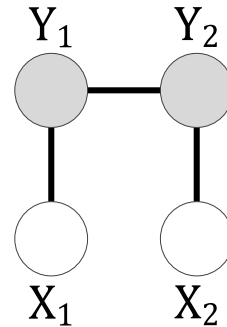
Conditional $P(Y_2 X_2, Y_1)$		$Y_2$	
		3	4
$Y_1=1$	$X_2=I$	$1-\epsilon$	$\epsilon$
	$X_2=O$	0.5	0.5
$Y_1=2$	$X_2=I$	0.5	0.5
	$X_2=O$	$\epsilon$	$1-\epsilon$

$Y_1, Y_2$	$P(Y_1 R)$	$P(Y_2 Y_1, O)$	$P(Y_1, Y_2 RO)$
1,3	0.8	0.5	0.4
1,4	0.8	0.5	0.4
2,3	0.2	$\epsilon$	$\epsilon$
2,4	0.2	$1-\epsilon$	0.2

## The Label Bias Problem: Potentials

Toy Problem – fragment of a CRF

$$\Psi(X, Y_1, Y_2) = \Psi(X_1, Y_1)\Psi(Y_1, Y_2)\Psi(X_2, Y_2)$$



$$Y_1 \in \{1,2\} \quad Y_2 \in \{3,4\}$$

Training Data:

$$8: X = \{RI\} \quad Y = \{13\}$$

$$2: X = \{RO\} \quad Y = \{24\}$$

# The Label Bias Problem: Potentials

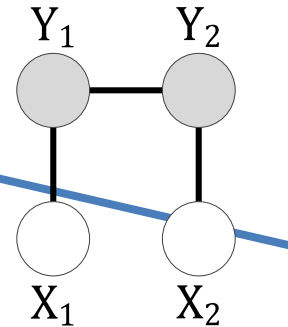
Toy Problem – fragment of a CRF

$\Psi$		$Y_1$	
		1	2
$X_1$	R	8	2
	-	$\epsilon$	$\epsilon$

$\Psi$		$Y_2$	
		3	4
$Y_1$	1	8	$\epsilon$
	2	$\epsilon$	2

$\Psi$		$Y_2$	
		3	4
$X_2$	I	8	$\epsilon$
	O	$\epsilon$	2

$$\Psi(X, Y_1, Y_2) = \Psi(X_1, Y_1)\Psi(Y_1, Y_2)\Psi(X_2, Y_2)$$



$Y_1 \in \{1,2\} Y_2 \in \{3,4\}$

Training Data:

8:  $X = \{RI\} Y = \{13\}$

2:  $X = \{RO\} Y = \{24\}$

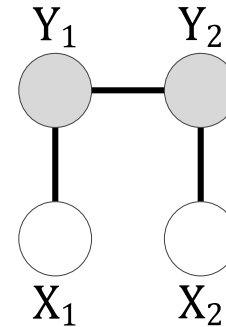
# The Label Bias Problem: Potentials

Toy Problem – fragment of a CRF

$\Psi$		$Y_1$	
		1	2
$X_1$	R	8	2
	-	$\epsilon$	$\epsilon$

$\Psi$		$Y_2$	
		3	4
$Y_1$	1	8	$\epsilon$
	2	$\epsilon$	2

$\Psi$		$Y_2$	
		3	4
$X_2$	I	8	$\epsilon$
	O	$\epsilon$	2



$Y_1 \in \{1,2\} Y_2 \in \{3,4\}$

Training Data:

8:  $X = \{RI\} Y = \{13\}$

2:  $X = \{RO\} Y = \{24\}$

$$\Psi(X, Y_1, Y_2) = \Psi(X_1, Y_1) \Psi(Y_1, Y_2) \Psi(X_2, Y_2)$$

$Y_1, Y_2$	$\Psi(X_1, Y_1)$	$\Psi(Y_1, Y_2)$	$\Psi(X_2, Y_2)$	$\Psi(Y_1, Y_2, X=RO)$	$P(Y_1, Y_2   X)$
1,3	8	8	$\epsilon$		
1,4	8	$\epsilon$	2		
2,3	2	$\epsilon$	$\epsilon$		
2,4	2	2	2		

Which labeling wins, now?



## The Label Bias Problem: Potentials

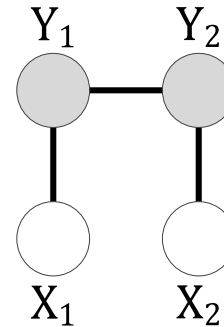
Toy Problem – fragment of a CRF

$\Psi$		$Y_1$	
		1	2
$X_1$	R	8	2
	-	$\epsilon$	$\epsilon$

$\Psi$		$Y_2$	
		3	4
$Y_1$	1	8	$\epsilon$
	2	$\epsilon$	2

$\Psi$		$Y_2$	
		3	4
$X_2$	I	8	$\epsilon$
	O	$\epsilon$	2

$$\Psi(X, Y_1, Y_2) = \Psi(X_1, Y_1)\Psi(Y_1, Y_2)\Psi(X_2, Y_2)$$



$$Y_1 \in \{1,2\} \quad Y_2 \in \{3,4\}$$

Training Data:

$$8: X = \{RI\} \quad Y = \{13\}$$

$$2: X = \{RO\} \quad Y = \{24\}$$

$Y_1, Y_2$	$\Psi(X_1, Y_1)$	$\Psi(Y_1, Y_2)$	$\Psi(X_2, Y_2)$	$\Psi(Y_1, Y_2, X=RO)$	$P(Y_1, Y_2 X)$
1,3	8	8	$\epsilon$	$64\epsilon$	small
1,4	8	$\epsilon$	2	$16\epsilon$	tiny
2,3	2	$\epsilon$	$\epsilon$	$2\epsilon^2$	infinitesimal
2,4	2	2	2	8	~100%

Because potentials do not have to normalize into probabilities until AFTER aggregation, they don't suffer from inappropriate conditioning.

Fun fact: We have seen this in class before!

- Graphical model  $\mathbf{G}(\mathbf{E}, \mathbf{V})$  with joint potential  $\Psi(\mathbf{V})$ ,  $\mathcal{C}$  the set of *cliques* in  $\mathbf{G}$  with  $c \in \mathcal{C}$  having edges  $\mathbf{E}_c$  and vertices  $\mathbf{V}_c$

$$P(\mathbf{V}) \propto \Psi(\mathbf{V}) = \prod_{c \in \mathcal{C}} \Psi_c(\mathbf{V}_c)$$

## Fun fact: We have seen this in class before!

- Graphical model  $\mathbf{G}(\mathbf{E}, \mathbf{V})$  with joint potential  $\Psi(\mathbf{V})$ ,  $\mathcal{C}$  the set of *cliques* in  $\mathbf{G}$  with  $c \in \mathcal{C}$  having edges  $\mathbf{E}_c$  and vertices  $\mathbf{V}_c$

$$P(\mathbf{V}) \propto \Psi(\mathbf{V}) = \prod_{c \in \mathcal{C}} \Psi_c(\mathbf{V}_c)$$

- $\mathbf{M}^3$  nets: cliques are pairs, and all conditioned on observed  $\mathbf{x}$ :

$$P(\mathbf{y}|\mathbf{x}) \propto \prod_{(i,j) \in \mathbf{E}} \psi_{ij}(y_i, y_j, \mathbf{x})$$

Fun fact: We have seen this in class before!

- Graphical model  $\mathbf{G}(\mathbf{E}, \mathbf{V})$  with joint potential  $\Psi(\mathbf{V})$ ,  $\mathcal{C}$  the set of *cliques* in  $\mathbf{G}$  with  $c \in \mathcal{C}$  having edges  $\mathbf{E}_c$  and vertices  $\mathbf{V}_c$

$$P(\mathbf{V}) \propto \Psi(\mathbf{V}) = \prod_{c \in \mathcal{C}} \Psi_c(\mathbf{V}_c)$$

- **M<sup>3</sup> nets**: cliques are pairs, and all conditioned on observed  $\mathbf{x}$ :

$$P(\mathbf{y}|\mathbf{x}) \propto \prod_{(i,j) \in \mathbf{E}} \psi_{ij}(y_i, y_j, \mathbf{x})$$

- **AMN**: cliques are pairs of nodes and singletons:

$$P_\phi(y) = \frac{1}{Z} \prod_i^N \phi_i(y_i) \prod_{i,j \in \mathbf{E}} \phi_{i,j}(y_i, y_j)$$

Where do Parameters Come From?

CRF's are part of the same general class,  $P(\mathbf{V}) \propto \Psi(\mathbf{V}) = \prod_{c \in \mathcal{C}} \Psi_c(\mathbf{V}_c)$

Where do Parameters Come From?

CRF's are part of the same general class,  $P(\mathbf{V}) \propto \Psi(\mathbf{V}) = \prod_{c \in \mathcal{C}} \Psi_c(\mathbf{V}_c)$

For trees, cliques are pairs of vertices sharing an edge ( $\mathbf{y}|_e$ ), and single vertices ( $\mathbf{y}|_v$ ):

$$\Psi(\mathbf{V}) = \prod_{e \in \text{Edge}} \Psi_e(\mathbf{y}|_e) \prod_{v \in \mathbf{V}} \Psi_v(\mathbf{y}|_v)$$

Where do Parameters Come From?

CRF's are part of the same general class,  $P(\mathbf{V}) \propto \Psi(\mathbf{V}) = \prod_{c \in \mathcal{C}} \Psi_c(\mathbf{V}_c)$

For trees, cliques are pairs of vertices sharing an edge ( $\mathbf{y}|_e$ ), and single vertices ( $\mathbf{y}|_v$ ):

$$\Psi(\mathbf{V}) = \prod_{e \in \text{Edge}} \Psi_e(\mathbf{y}|_e) \prod_{v \in \mathbf{V}} \Psi_v(\mathbf{y}|_v)$$

And because this is a conditional network with  $\mathbf{V} = \mathbf{X} \cup \mathbf{Y}$

$$\Psi(\mathbf{Y}|\mathbf{X}) = \prod_{e \in \text{Edge}} \Psi_e(\mathbf{y}|_e, x) \prod_{v \in \mathbf{V}} \Psi_v(\mathbf{y}|_v, x)$$

## Where do Parameters Come From?

CRF's are part of the same general class,  $P(\mathbf{V}) \propto \Psi(\mathbf{V}) = \prod_{c \in \mathcal{C}} \Psi_c(\mathbf{V}_c)$

For trees, cliques are pairs of vertices sharing an edge ( $\mathbf{y}|_e$ ), and single vertices ( $\mathbf{y}|_v$ ):

$$\Psi(\mathbf{V}) = \prod_{e \in \text{Edge}} \Psi_e(\mathbf{y}|_e) \prod_{v \in \mathbf{V}} \Psi_v(\mathbf{y}|_v)$$

And because this is a conditional network with  $\mathbf{V} = \mathbf{X} \cup \mathbf{Y}$

$$\Psi(\mathbf{Y}|\mathbf{X}) = \prod_{e \in \text{Edge}} \Psi_e(\mathbf{y}|_e, x) \prod_{v \in \mathbf{V}} \Psi_v(\mathbf{y}|_v, x)$$

An exponential identity gives us

$$\Psi(\mathbf{Y}|\mathbf{X}) = \exp\left(\sum_{e \in \mathbf{E}} \log\left(\Psi_e(\mathbf{y}|_e, x)\right) + \sum_{v \in \mathbf{V}} \log\left(\Psi_v(\mathbf{y}|_v, x)\right)\right)$$



## Where do Parameters Come From?

CRF's are part of the same general class,  $P(\mathbf{V}) \propto \Psi(\mathbf{V}) = \prod_{c \in \mathcal{C}} \Psi_c(\mathbf{V}_c)$

For trees, cliques are pairs of vertices sharing an edge ( $\mathbf{y}|_e$ ), and single vertices ( $\mathbf{y}|_v$ ):

$$\Psi(\mathbf{V}) = \prod_{e \in \text{Edge}} \Psi_e(\mathbf{y}|_e) \prod_{v \in \mathbf{V}} \Psi_v(\mathbf{y}|_v)$$

And because this is a conditional network with  $\mathbf{V} = \mathbf{X} \cup \mathbf{Y}$

$$\Psi(\mathbf{Y}|\mathbf{X}) = \prod_{e \in \text{Edge}} \Psi_e(\mathbf{y}|_e, x) \prod_{v \in \mathbf{V}} \Psi_v(\mathbf{y}|_v, x)$$

An exponential identity gives us

$$\Psi(\mathbf{Y}|\mathbf{X}) = \exp\left(\sum_{e \in \mathbf{E}} \log\left(\Psi_e(\mathbf{y}|_e, x)\right) + \sum_{v \in \mathbf{V}} \log\left(\Psi_v(\mathbf{y}|_v, x)\right)\right)$$

Potentials can be ANY positive values... like linear combinations of arbitrary features

$$\Psi(\mathbf{Y}|\mathbf{X}) = \exp\left(\sum_{e \in \mathbf{E}, k \in \mathbf{K}} \lambda_k f_k(\mathbf{e}, \mathbf{y}|_e, x) + \sum_{v \in \mathbf{V}, k' \in \mathbf{K}'} \mu_{k'} g_{k'}(\mathbf{v}, \mathbf{y}|_v, x)\right)$$

# Improved iterative scaling

- Want to maximize log-likelihood with respect to parameters

$$\theta = (\lambda_1, \lambda_2, \dots; \mu_1, \mu_2, \dots)$$

---

<sup>1</sup>Della Pietra *et al.* (1997)

## Improved iterative scaling

- Want to maximize log-likelihood with respect to parameters

$$\theta = (\lambda_1, \lambda_2, \dots; \mu_1, \mu_2, \dots)$$

- Method: *Improved Iterative Scaling*<sup>1</sup>: Extension of Generalised Iterative Scaling (Darroch and Ratcliff 1972).

---

<sup>1</sup>Della Pietra *et al.* (1997)

## Improved iterative scaling

- Want to maximize log-likelihood with respect to parameters

$$\theta = (\lambda_1, \lambda_2, \dots; \mu_1, \mu_2, \dots)$$

- Method: *Improved Iterative Scaling*<sup>1</sup>: Extension of Generalised Iterative Scaling (Darroch and Ratcliff 1972).
- *Improved* because features need not sum to constant.

---

<sup>1</sup>Della Pietra *et al.* (1997)

## Improved iterative scaling

- Want to maximize log-likelihood with respect to parameters

$$\theta = (\lambda_1, \lambda_2, \dots; \mu_1, \mu_2, \dots)$$

- Method: *Improved Iterative Scaling*<sup>1</sup>: Extension of Generalised Iterative Scaling (Darroch and Ratcliff 1972).
- *Improved* because features need not sum to constant.
- Idea: new set of parameters

$\theta' = \theta + \delta\theta = (\lambda_1 + \delta\lambda_1, \dots; \mu_1 + \delta\mu_1, \dots)$  which will not decrease objective function. Iteratively apply!

---

<sup>1</sup>Della Pietra *et al.* (1997)

## Improved iterative scaling

- Want to maximize log-likelihood with respect to parameters

$$\theta = (\lambda_1, \lambda_2, \dots; \mu_1, \mu_2, \dots)$$

- Method: *Improved Iterative Scaling*<sup>1</sup>: Extension of Generalised Iterative Scaling (Darroch and Ratcliff 1972).
- *Improved* because features need not sum to constant.
- Idea: new set of parameters

$$\theta' = \theta + \delta\theta = (\lambda_1 + \delta\lambda_1, \dots; \mu_1 + \delta\mu_1, \dots)$$

which will not decrease objective function. Iteratively apply!

- Problem: slow, and nobody uses this any more.

<sup>1</sup>Della Pietra *et al.* (1997)

# Modern CRF training - L-BFGS

- Generally use L-BFGS<sup>2</sup> algorithm.

---

<sup>2</sup>Limited-Memory Broyden-Fletcher-Goldfarb-Shanno Algorithm

## Modern CRF training - L-BFGS

- Generally use L-BFGS<sup>2</sup> algorithm.
- Approximates Newton's method. Optimise multivariate function  $f(\theta)$  through updates

$$\theta_{t+1} = \theta_t - [Hf(\theta_t)]^{-1} \nabla f(\theta_t)$$

---

<sup>2</sup>Limited-Memory Broyden-Fletcher-Goldfarb-Shanno Algorithm



## Modern CRF training - L-BFGS

- Generally use L-BFGS<sup>2</sup> algorithm.
- Approximates Newton's method. Optimise multivariate function  $f(\theta)$  through updates

$$\theta_{t+1} = \theta_t - [Hf(\theta_t)]^{-1} \nabla f(\theta_t)$$

- *Quasi-Newtonian*: approximates Hessian  $Hf(\theta)$ .

---

<sup>2</sup>Limited-Memory Broyden-Fletcher-Goldfarb-Shanno Algorithm

## Modern CRF training - L-BFGS

- Generally use L-BFGS<sup>2</sup> algorithm.
- Approximates Newton's method. Optimise multivariate function  $f(\theta)$  through updates

$$\theta_{t+1} = \theta_t - [Hf(\theta_t)]^{-1} \nabla f(\theta_t)$$

- *Quasi-Newtonian*: approximates Hessian  $Hf(\theta)$ .
- Limited-memory: doesn't store full (approximate) Hessian.

---

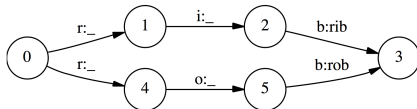
<sup>2</sup>Limited-Memory Broyden-Fletcher-Goldfarb-Shanno Algorithm 

## Label bias

- Generate data with noisy HMM.
- 4-state system (not counting 'initial state'), transitions:

- $1 \Rightarrow 2 \Rightarrow 3$

- $4 \Rightarrow 5 \Rightarrow 3$



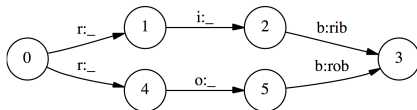
- Emissions: highly biased!
  - $P(X = Y\text{'s preferred value} | Y) = 29/32$
  - $P(X = \text{other} | Y) = 1/32$
- Preferred values:  $1 \rightarrow \text{'r'}$ ,  $4 \rightarrow \text{'r'}$ ,  $2 \rightarrow \text{'i'}$ ,  $5 \rightarrow \text{'o'}$ ,  $3 \rightarrow \text{'b'}$ .

## Label bias

- Generate data with noisy HMM.
- 4-state system (not counting 'initial state'), transitions:

- $1 \Rightarrow 2 \Rightarrow 3$

- $4 \Rightarrow 5 \Rightarrow 3$

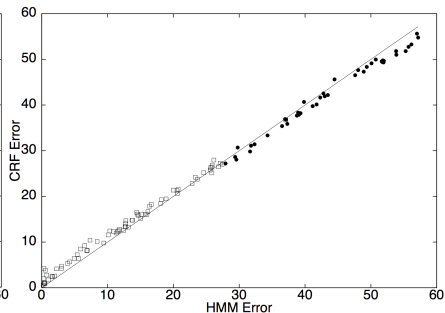
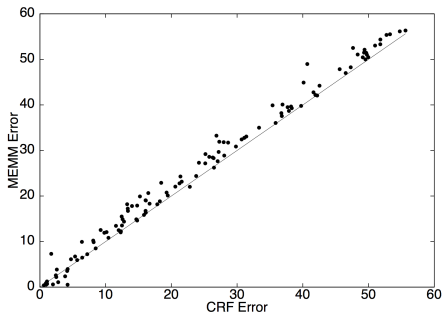


- Emissions: highly biased!
  - $P(X = Y\text{'s preferred value} | Y) = 29/32$
  - $P(X = \text{other} | Y) = 1/32$
- Preferred values:  $1 \rightarrow \text{'r'}$ ,  $4 \rightarrow \text{'r'}$ ,  $2 \rightarrow \text{'i'}$ ,  $5 \rightarrow \text{'o'}$ ,  $3 \rightarrow \text{'b'}$ .
- Result: CRF error 4.6%, MEMM error 42%.

## Mixed-order sources

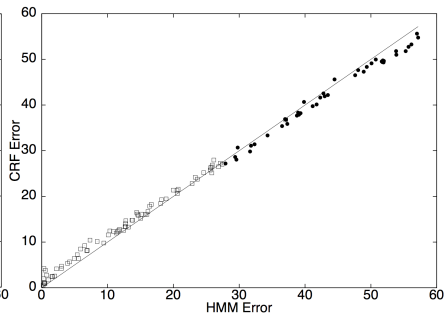
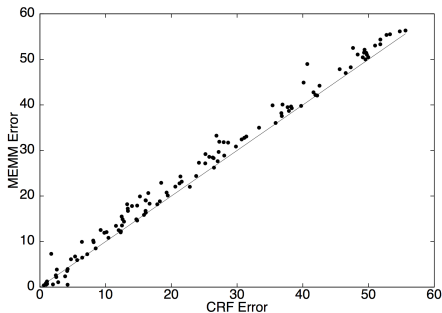
- Generate data with mixed-order HMM:
  - Transitions:  $(1 - \alpha)p_1(\mathbf{y}_i|\mathbf{y}_{i-1}) + \alpha p_2(\mathbf{y}_i|\mathbf{y}_{i-1}, \mathbf{y}_{i-2})$
  - Emissions:  $(1 - \alpha)p_1(\mathbf{x}_i|\mathbf{y}_i) + \alpha p_2(\mathbf{x}_i|\mathbf{y}_i, \mathbf{x}_{i-1})$
- Five labels, 26 observation values.
- Training/testing: 1000 sequences of length 25.
- CRF trained with Algorithm S (modified IIS). MEMM trained with iterative scaling.
- Viterbi to label test set.

## Mixed-order sources: results



■ Squares :  $\alpha < 0.5$ .

## Mixed-order sources: results



■ Squares :  $\alpha < 0.5$ .

■ CRF sort of wins?

# Part of Speech Tagging

- Penn Treebank: 45 syntactic tags, label each word in sentence.
- Train first-order HMM, MEMM, CRF.



## Part of Speech Tagging

- Penn Treebank: 45 syntactic tags, label each word in sentence.
- Train first-order HMM, MEMM, CRF.

<i>model</i>	<i>error</i>	<i>oov error</i>
HMM	5.69%	45.99%
MEMM	6.37%	54.61%
CRF	5.55%	48.05%
MEMM <sup>+</sup>	4.81%	26.99%
CRF <sup>+</sup>	4.27%	23.76%

<sup>+</sup>Using spelling features

## Part of Speech Tagging

- Penn Treebank: 45 syntactic tags, label each word in sentence.
- Train first-order HMM, MEMM, CRF.

<i>model</i>	<i>error</i>	<i>oov error</i>
HMM	5.69%	45.99%
MEMM	6.37%	54.61%
CRF	5.55%	48.05%
MEMM <sup>+</sup>	4.81%	26.99%
CRF <sup>+</sup>	4.27%	23.76%

<sup>+</sup>Using spelling features

- Spelling features exploit conditional framework.
- Examples: starts with number/upper case?, contains hyphen, has suffix?

# Skip-chain CRF

- Example: skip-chain CRF<sup>3</sup>.

---

<sup>3</sup>From *An Introduction to Conditional Random Fields for Relational*

*Learning*, Charles Sutton and Andrew McCallum, 2006

# Skip-chain CRF

- Example: skip-chain CRF<sup>3</sup>.
- Has long-range features!

---

<sup>3</sup>From *An Introduction to Conditional Random Fields for Relational*

*Learning*, Charles Sutton and Andrew McCallum, 2006

# Skip-chain CRF

- Example: skip-chain CRF<sup>3</sup>.
- Has long-range features!
- Basic idea: extend linear-chain CRF by joining some distant observations with ‘skip edges’.

---

<sup>3</sup>From *An Introduction to Conditional Random Fields for Relational*

*Learning*, Charles Sutton and Andrew McCallum, 2006

# Skip-chain CRF

- Example: skip-chain CRF<sup>3</sup>.
- Has long-range features!
- Basic idea: extend linear-chain CRF by joining some distant observations with ‘skip edges’.
- Connect multiple mentions of entity across whole document.

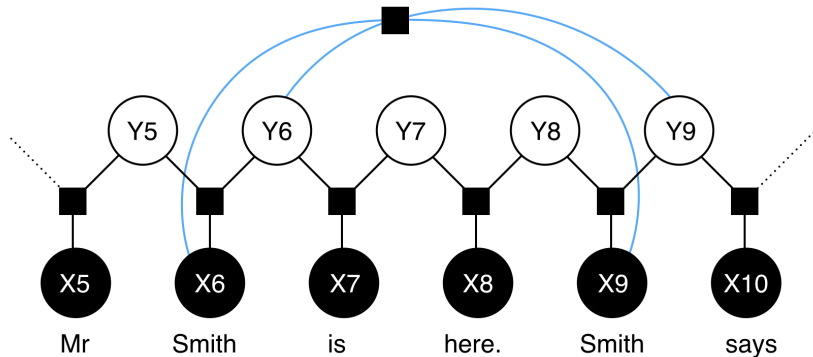
---

<sup>3</sup>From *An Introduction to Conditional Random Fields for Relational*

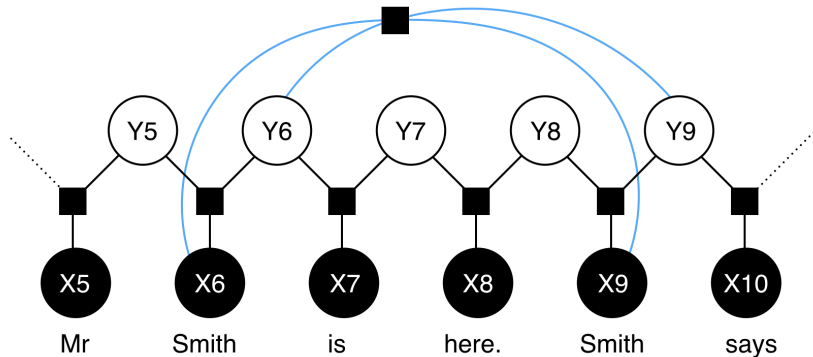
*Learning*, Charles Sutton and Andrew McCallum, 2006

## Skip-chain CRF

**Example:** Note: Squares denote factors (e.g. potential functions).



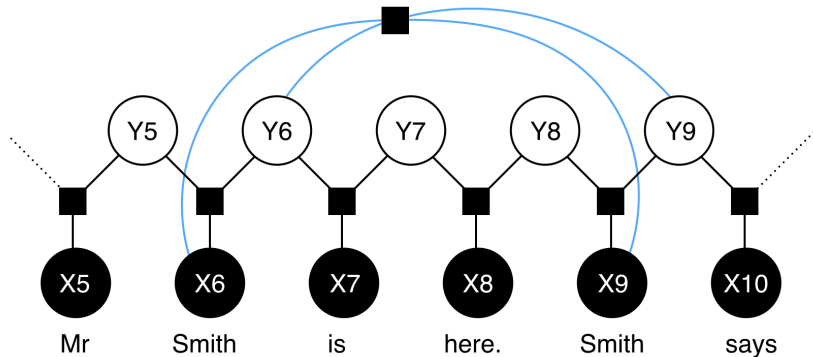
Example: Note: Squares denote factors (e.g. potential functions).



Question: Ignoring the skip edges (in blue), what potentials does  $Y_i$  appear in?



Example: Note: Squares denote factors (e.g. potential functions).



Question: Ignoring the skip edges (in blue), what potentials does

$Y_i$  appear in? Answer:  $\psi(Y_i, Y_{i-1}, X_i), \psi(Y_{i+1}, Y_i, X_{i+1})$

# Skip-chain task

- Data: 485 email announcements for seminars at CMU.

# Skip-chain task

- Data: 485 email announcements for seminars at CMU.
- Task: identify start time, end time, location, speaker.

# Skip-chain task

- Data: 485 email announcements for seminars at CMU.
- Task: identify start time, end time, location, speaker.
- Linear chain CRF with skip edges between identical capitalised words.

# Skip-chain task

- Data: 485 email announcements for seminars at CMU.
- Task: identify start time, end time, location, speaker.
- Linear chain CRF with skip edges between identical capitalised words.
- Other word-specific features e.g. ‘appears in list of first names’, ‘upper case’, ‘appears to be part of time/date’ (by regex), etc.

## Skip-chain CRF

## Skip-chain results

System	stime	etime	location	speaker	overall
BIEN Peshkin and Pfeffer [2003]	96.0	<b>98.8</b>	87.1	76.9	89.7
Linear-chain CRF	<b>97.5</b>	97.5	<b>88.3</b>	77.3	90.2
Skip-chain CRF	96.7	97.2	88.1	<b>80.4</b>	<b>90.6</b>

- Values are F1 scores.

## Skip-chain results

System	stime	etime	location	speaker	overall
BIEN Peshkin and Pfeffer [2003]	96.0	<b>98.8</b>	87.1	76.9	89.7
Linear-chain CRF	<b>97.5</b>	97.5	<b>88.3</b>	77.3	90.2
Skip-chain CRF	96.7	97.2	88.1	<b>80.4</b>	<b>90.6</b>

- Values are F1 scores.
- Repeated occurrences of speaker improve skip-chain performance.

## Skip-chain results

System	stime	etime	location	speaker	overall
BIEN Peshkin and Pfeffer [2003]	96.0	<b>98.8</b>	87.1	76.9	89.7
Linear-chain CRF	<b>97.5</b>	97.5	<b>88.3</b>	77.3	90.2
Skip-chain CRF	96.7	97.2	88.1	<b>80.4</b>	<b>90.6</b>

- Values are F1 scores.
- Repeated occurrences of speaker improve skip-chain performance.
- Tokens are *consistently* classified by skip-chain. Linear-chain is inconsistent on **30.2** speakers, skip-chain: **4.8**.



# Summary

- CRFs combine *discriminative* (e.g. MEMM) and *undirected* (e.g. MRF) properties to solve problems:

## Summary

- CRFs combine *discriminative* (e.g. MEMM) and *undirected* (e.g. MRF) properties to solve problems:
  - Global normalisation avoids label bias.

# Summary

- CRFs combine *discriminative* (e.g. MEMM) and *undirected* (e.g. MRF) properties to solve problems:
  - Global normalisation avoids label bias.
  - Conditioning on observations avoids modelling complex dependencies.

## Summary

- CRFs combine *discriminative* (e.g. MEMM) and *undirected* (e.g. MRF) properties to solve problems:
  - Global normalisation avoids label bias.
  - Conditioning on observations avoids modelling complex dependencies.
  - Enables use of features using global structure.

## Summary

- CRFs combine *discriminative* (e.g. MEMM) and *undirected* (e.g. MRF) properties to solve problems:
  - Global normalisation avoids label bias.
  - Conditioning on observations avoids modelling complex dependencies.
  - Enables use of features using global structure.
- Examples in paper strangely insubstantial, but CRFs are widely and successfully used.